

A Lower Bound for Succinct Rank Queries

Mihai Pătraşcu
IBM Almaden

July 6, 2009

Abstract

The rank problem in succinct data structures asks to preprocess an array $A[1..n]$ of bits into a data structure using as close to n bits as possible, and answer queries of the form $\text{RANK}(k) = \sum_{i=1}^k A[i]$. The problem has been intensely studied, and features as a subroutine in a majority of succinct data structures.

We show that in the cell probe model with w -bit cells, if rank takes t time, the space of the data structure must be at least $n + n/w^{O(t)}$ bits. This redundancy/query trade-off is essentially optimal, matching our upper bound from [FOCS'08].

1 Introduction

1.1 The Complexity of Rank

Consider an array $A[1..n]$ of bits. Can we preprocess this array into a data structure of size $n + r$ bits, for small redundancy r , which supports rank queries $\text{RANK}(k) = \sum_{i=1}^k A[i]$ efficiently? The problem of supporting rank (and the related select queries) is the bread-and-butter of succinct data structures. It finds use in most other data structures (for representing trees, graphs, suffix trees / suffix arrays etc), and its redundancy / query trade-off has come under quite a bit of attention.

Rank already had a central position in the seminal papers on succinct data structures. Jacobson [Jac89], in FOCS'89, and Clark and Munro [CM96], in SODA'96, gave the first data structures using space $n + o(n)$ and constant query time. These results were slightly improved in [Mun96, MRR01, RRR02].

In several applications, the set of ones is not dense in the array. Thus, the problem was generalized to storing an array $A[1..u]$, containing n ones and $u - n$ zeros. The optimal space is $B = \lg \binom{u}{n}$. Pagh [Pag01] achieved space $B + O(n \cdot \frac{(\lg \lg n)^2}{\lg n})$ for this sparse problem. Recently, Golynski et al. [GGG⁺07] achieved $B + O(n \cdot \frac{\lg \lg u}{\lg^2 n})$. Subsequently, Golynski et al. [GRR08] have achieved space $B + O(n \cdot \frac{\lg \lg n \cdot \lg(u/n)}{\lg^2 n})$.

In my paper from FOCS'08 [Păt08], I gave a qualitative improvement to these bounds, showing an exponential dependence between the query time and the redundancy. Specifically, with query time $O(t)$, the achievable redundancy is $r \leq n / (\frac{\lg n}{t})^t$. This improved the redundancy for many succinct data structures where rank/select queries were the bottleneck.

Given the surprising nature of this improvement, a natural question is whether we can do much better. In this paper, we show that we cannot, at least for the basic rank queries:

Theorem 1. *In the cell-probe model with words of $w \geq \lg n$ bits, a data structure that supports rank queries in t cell probes requires at least $n + n/w^{O(t)}$ bits of space.*

All succinct data structure papers assume $w = \lg n$. The lower bound matches my upper bound, except for the difference between $(\lg n)^t$ and $(\frac{\lg n}{t})^t$. This difference is inconsequential for small $t < \lg^{0.99} n$. If we want a polynomially small redundancy (say, less than n^α , for some constant $\alpha < 1$), the upper bound says that $t = O(\lg n)$ is sufficient. The lower bound says that $t = \Omega(\lg n / \lg \lg n)$ is necessary. It is unclear which bound is the optimal one in this regime.

1.2 Lower Bounds for Succinct Data Structures

Much work in lower bounds for succinct data structures has been in the so-called systematic model. In this model, the array A must be represented as is, i.e. the data structure only has oracle access to it (it can read any w consecutive bits at $O(1)$ cost). In addition, the data structure may store an index of sublinear size, which the query algorithm can examine at no cost. See [GM03, Mil05, GRR08, Gol07] for increasingly tight lower bounds in this model. Note, however, that in the systematic model, the best achievable redundancy with query time t is $\frac{n}{t \cdot \text{poly} \lg n}$, i.e. there is a linear trade-off between redundancy and query time. This is significantly improved by my (non-systematic) upper bounds [Pät08], and these lower bounds qualitatively miss the nature of this improvement.

In the unrestricted cell-probe model, the first lower bounds were shown by Gál and Miltersen [GM03] in 2003. These lower bounds were strong, showing a linear dependence between the query and redundancy $r \cdot t = \Omega(n / \lg n)$. However, the problem being analyzed is somewhat unnatural: the bound applies to polynomial evaluation, for which nontrivial succinct upper bounds appear unlikely. Their technique, which is based on the strong error correction implicit in their problem, remains powerless for “easier” problems. (Thus, succinct data structures are unusual for lower bounds, in that the difficult goal seems to be proving *lower* lower bounds for natural problems.)

A significant break-through occurred in SODA’09, when Golynski [Gol09] showed a lower bound of $r \cdot t^2 = \Omega(n)$ for the problem of storing a permutation and querying $\pi(\cdot)$ and $\pi^{-1}(\cdot)$. This quadratic trade-off is tight for storing a permutation and its inverse. Golynski’s technique is based on the inherent difficulty of storing a function and its inverse without doubling the space. However, due to the particular attention it pays to inverses, it is unclear how it could generalize to problems like rank.

In this paper, we make further progress on getting lower bounds for natural problems, and analyze one of the central problems in succinct data structures. It is reasonable to hope that our lower bound technique will generalize to many other problems, given the many applications of rank queries.

2 The Proof

2.1 An Entropy Bound

The structure of the rank problem is not particularly important in the lower bound proof. All that is needed is an inequality on the entropy of rank queries that we describe here. Essentially, the lower bound applies to any problem which satisfies a similar entropy condition.

The possible queries come from the universe $[n]$. Imagine that this universe is divided into k blocks of equal size (the remainder is ignored if k doesn't divide n). Let $Q_\Delta \subset [n]$ be the set containing the Δ -th query (counting from zero) in each block. For a set Q of queries, let $\text{Ans}(Q)$ be the vector of answers to the queries in Q . We treat $\text{Ans}(Q)$ as a random variable, depending on the random choice of the input $A[1..n]$.

Lemma 2. *Let A be chosen uniformly at random in $\{0,1\}^n$, and let Δ and any $Q^* \subseteq Q_\Delta$ be arbitrary. Then, for any event \mathcal{E} with $\Pr[\mathcal{E}] = 2^{-\varepsilon|Q^*|}$ for a small enough constant ε , we have:*

$$H(\text{Ans}(Q_0) \mid \mathcal{E}) + H(\text{Ans}(Q^*) \mid \mathcal{E}) - H(\text{Ans}(Q_0), \text{Ans}(Q^*) \mid \mathcal{E}) = \Omega(|Q^*|)$$

Proof. Let us ignore the conditioning on \mathcal{E} for now. The lemma says that representing the answers to the queries Q_0 and (a subset of) Q_Δ separately loses $\Omega(1)$ bits of entropy per block compared to the optimal joint encoding.

Let h_m be entropy of the binomial distribution on m unbiased trials. The entropy $H(\text{Ans}(Q_0))$ is exactly equal to $k \cdot h_{n/k}$: the answer of a query minus the answer of the previous is exactly a binomial on n/k random bits. In all blocks that do not contain an element of Q^* , the contribution of the block in $H(\text{Ans}(Q_0))$ is cancelled by its contribution in $H(\text{Ans}(Q_0), \text{Ans}(Q^*))$.

Blocks that contain an element from Q^* (except the first block) contribute:

- $h_{n/k}$ to $H(\text{Ans}(Q_0))$;
- at least $h_{n/k}$ to $H(\text{Ans}(Q^*))$. The contribution is more if the previous block did not contain an element from Q^* ;
- exactly $h_\Delta + h_{n/k-\Delta}$ to $H(\text{Ans}(Q_0), \text{Ans}(Q^*))$.

Thus, the block contributes $2h_{n/k} - h_\Delta - h_{n/k-\Delta}$ to the sum. Using the known estimation $h_m = \frac{1}{2} \ln(\frac{\pi e}{2} m) + O(\frac{1}{m})$, this quantity is minimized when $\Delta = \frac{n}{2k}$, and is always at least $\ln 2 - o(1)$.

The fact that conditioning on \mathcal{E} does not change the result comes from a standard independence trick in lower bounds. We decomposed $H(\text{Ans}(Q_0)) + H(\text{Ans}(Q^*)) - H(\text{Ans}(Q_0), \text{Ans}(Q^*))$ as the sum over Q^* independent variables (essentially¹). Each component was $\Omega(1)$ with constant probability. By a Chernoff bound, the sum is $\Omega(|Q^*|)$ with probability $2^{-\Omega(|Q^*|)}$. Thus, even if we condition on an event of probability $2^{-\varepsilon|Q^*|}$, the sum must remain $\Omega(|Q^*|)$ with overwhelming probability. \square

2.2 Cell-Probe Elimination

To support the induction in our proof, we augment the cell-probe model with *published bits*. These bits represent a memory of bounded size which the query algorithm can examine *at no cost*. Like the regular memory (which must be examined through cell probes), the published bits are initialized at construction time, as a function of the input $A[1..n]$. Observe that if we have n published bits, the problem can be solved trivially.

Our proof will try to publish a small number of cells from the regular memory which are accessed frequently. Thus, the complexity of many queries will decrease by at least one. The argument is

¹The careful reader has probably noticed that we actually decomposed it into *two* sums, each of which has Q^* terms independent among themselves; however, the sums are dependent. We are subtracting the entropy of sub-blocks of size Δ from the entropy of blocks of size n/k in the first sum; and the entropy of sub-blocks of size $n/k - \Delta$ from the entropy of blocks of size n/k in the second sum. The analysis proceeds by union bound over the two sums.

then applied iteratively: the cell-probe complexity decreases, as more and more bits are published. If we arrive at zero cell probes and less than n published bits, we have a contradiction.

Let $\text{Probes}(q)$ be the set of cells probed by query q ; this is a random variable, since the query can be adaptive. Also let $\text{Probes}(Q) = \bigcup_{q \in Q} \text{Probes}(q)$.

The main technical result in our proof is captured in the following lemma, the proof of which appears in the next section:

Lemma 3. *Assume a data structure uses $P = o(n)$ published bits, and at most n memory bits. Break the queries into $k = \gamma \cdot P$ blocks, for a large enough constant γ . Then:*

$$\Pr_{A, q \in [n]} [\text{Probes}(q) \cap \text{Probes}(Q_0) \neq \emptyset] = \Omega(1)$$

The lemma shows that $\text{Probes}(Q_0)$ are a good set of cells to publish, since a constant fraction of the queries probe at least one cell from this set.

Completing the proof is now easy. If the data structure has redundancy r , begin by publishing some arbitrary $P_0 = r$ bits, to satisfy the condition that there are at most n bits in regular memory.

In step $i = 0, 1, 2, \dots$, we let $k_i = \gamma \cdot P_i$, and publish the cells in $\text{Probes}(Q_0)$, together with their address. The number of published bits increases to $P_{i+1} = k_i \cdot (w + O(\lg n)) = O(P_i w)$. The cell-probe complexity of an average query decreases by $\Omega(1)$.

Since the average case complexity cannot go below zero, the number of iterations that we are able to make must be $O(t)$. The only reason we may fail to make another iteration is a violation to the lemma's condition $P = o(n)$. Thus, $P_{O(t)} = \Omega(n)$, that is $r \cdot w^{O(t)} \geq n$. This is the desired trade-off.

2.3 An Encoding Argument

In this section, we prove Lemma 3. Our proof is an encoding argument: we show that, if the conclusion of the lemma failed, we could encode a uniformly random A using strictly less than n bits.

Let P and k be as in our lemma's statement, and assume $\Pr_{A, q \in [n]} [\text{Probes}(q) \cap \text{Probes}(Q_0) \neq \emptyset] \leq \varepsilon$, for a small enough constant ε . We thus know that a random query is very likely to probe cells not in $\text{Probes}(Q_0)$.

By averaging, there exists a $\Delta \in \{1, \dots, n/k\}$ such that $\Pr_{A, q \in Q_\Delta} [\text{Probes}(q) \cap \text{Probes}(Q_0) \neq \emptyset] \leq \varepsilon$. We are only going to concentrate on the queries in Q_Δ .

More specifically, we are going to concentrate on the queries that probe no cell from $\text{Probes}(Q_0)$: $Q^* = \{q \in Q_\Delta \mid \text{Probes}(q) \cap \text{Probes}(Q_0) = \emptyset\}$. Note that $\mathbf{E}_A[|Q^*|] \geq (1 - \varepsilon)k$.

Intuitively speaking, our contradiction is found as follows. The answers to queries Q_0 must be encoded in the cells $\text{Probes}(Q_0)$. The answers to queries Q^* must be encoded in the cells $\text{Probes}(Q^*)$, which, by definition, is disjoint from $\text{Probes}(Q_0)$. But the answers $\text{Ans}(Q_0)$ and $\text{Ans}(Q^*)$ are highly correlated (by Lemma 2). Thus, if the two answers are written in disjoint sets of cells, a lot of entropy is being wasted, which is impossible for a succinct data structure.

The footprint. We first formalize the intuitive notion of “the contents of cells $\text{Probes}(Q)$.” Define the footprint $\text{Foot}(Q)$ of a query set Q by the following algorithm. We assume the published bits are known in the course of the definition. Enumerate queries $q \in Q$ in increasing order. For each query, simulate its execution one cell probe at a time. If a cell has already been included in the

footprint, ignore it. Otherwise, append the contents (but not the address) of the new cell in the footprint. Observe that $\text{Foot}(Q)$ is a string of exactly $|\text{Probes}(Q)| \cdot w$ bits.

We observe that $\text{Ans}(Q)$ is a function of $\text{Foot}(Q)$ and the published bits. Indeed, we can simulate the queries in order. At each step, we know how the query algorithm acts based on the published bits and the previously read cells. Thus, we know the address of the next cell to be read. We can check whether the cell was already in the footprint (since we also know the address of previous cells). If not, we read the next w bits of the footprint, which are precisely the contents of this cell, and continue the simulation.

The encoding. Our encoding for the array A will consist of the following:

1. the published bits (P bits). Denote these bits by the random variable \mathcal{P} .
2. the identity of the set Q^* as a subset of Q_Δ . This uses $O(\lg \binom{k}{|Q^*|}) = O(\lg \binom{k}{k-|Q^*|})$ bits. By submodularity, the average length of this component is on the order of:

$$\mathbf{E}\left[\lg \binom{k}{k-|Q^*|}\right] \leq \lg \left(\mathbf{E}\left[\binom{k}{k-|Q^*|}\right]\right) \leq \lg \left(\frac{k}{\varepsilon k}\right) = k \cdot O(\varepsilon \lg \frac{1}{\varepsilon})$$

3. the answers $\text{Ans}(Q_0 \cup Q^*)$, encoded jointly. Using Huffman coding, this requires $H(\text{Ans}(Q_0 \cup Q^*)) + O(1)$ bits on average.
4. the footprint $\text{Foot}(Q_0)$, encoded optimally given the knowledge of $\text{Ans}(Q_0)$ and the published bits. This takes $H(\text{Foot}(Q_0) \mid \text{Ans}(Q_0), \mathcal{P}) + O(1)$ bits on average.
5. the footprint $\text{Foot}(Q^*)$, encoded optimally given the knowledge of Q^* , $\text{Ans}(Q^*)$, and the published bits. This takes $H(\text{Foot}(Q^*) \mid Q^*, \text{Ans}(Q^*), \mathcal{P}) + O(1)$ bits on average.
6. all cells outside $\text{Probes}(Q_0) \cup \text{Probes}(Q^*)$, included verbatim with w bits per cell. As noted above, the cell addresses $\text{Probes}(Q_0)$ and $\text{Probes}(Q^*)$ can be decoded from $\text{Foot}(Q_0)$, respectively $\text{Foot}(Q^*)$, and the published bits. Thus, we know exactly which cells to include in this component. This part takes $n - \mathbf{E}[|\text{Probes}(Q_0)| + |\text{Probes}(Q^*)|] \cdot w$ bits on average.

Observe that this encoding includes the published bits and all cells in the memory (though the cells in $\text{Probes}(Q_0)$ and $\text{Probes}(Q^*)$ are included in a compressed format). Thus, all n queries can be simulated. If all n answers are known, the array A can be decoded. Thus, this is a valid encoding of A .

It remains to analyze the average size of the encoding. To bound item 4., we can write:

$$H(\text{Foot}(Q_0) \mid \text{Ans}(Q_0), \mathcal{P}) = H(\text{Foot}(Q_0), \text{Ans}(Q_0), \mathcal{P}) - H(\text{Ans}(Q_0), \mathcal{P})$$

But $H(\text{Foot}(Q_0), \text{Ans}(Q_0), \mathcal{P}) = H(\text{Foot}(Q_0), \mathcal{P})$, since the answers can be decoded from the footprint and the published bits. Now note that $H(\text{Foot}(Q_0), \mathcal{P}) \leq \mathbf{E}[|\text{Probes}(Q_0)|] \cdot w + P$, since this is the size in bits of the footprint and the published bits. Finally, note that $H(\text{Ans}(Q_0), \mathcal{P}) \geq H(\text{Ans}(Q_0))$. Thus:

$$H(\text{Foot}(Q_0) \mid \text{Ans}(Q_0), \mathcal{P}) \leq \mathbf{E}[|\text{Probes}(Q_0)|] \cdot w + P - H(\text{Ans}(Q_0))$$

Similarly, item 5. is bounded by:

$$H(\text{Foot}(Q^*) \mid Q^*, \text{Ans}(Q^*), \mathcal{P}) \leq \mathbf{E}[|\text{Probes}(Q^*)|] \cdot w + P + k \cdot O(\varepsilon \lg \frac{1}{\varepsilon}) - H(Q^*, \text{Ans}(Q^*))$$

Summing up all components, our encoding has expected size:

$$n + 3P + k \cdot O(\varepsilon \lg \frac{1}{\varepsilon}) + H(\text{Ans}(Q_0), \text{Ans}(Q^*)) - H(\text{Ans}(Q_0)) - H(Q^*, \text{Ans}(Q^*)) \quad (1)$$

We can now rewrite:

$$\begin{aligned} & H(\text{Ans}(Q_0)) + H(Q^*, \text{Ans}(Q^*)) - H(\text{Ans}(Q_0), \text{Ans}(Q^*)) \\ \geq & H(\text{Ans}(Q_0) \mid Q^*) + H(\text{Ans}(Q^*) \mid Q^*) + H(Q^*) - H(\text{Ans}(Q_0), \text{Ans}(Q^*), Q^*) \\ = & H(\text{Ans}(Q_0) \mid Q^*) + H(\text{Ans}(Q^*) \mid Q^*) + H(Q^*) - H(\text{Ans}(Q_0), \text{Ans}(Q^*) \mid Q^*) - H(Q^*) \\ = & H(\text{Ans}(Q_0) \mid Q^*) + H(\text{Ans}(Q^*) \mid Q^*) - H(\text{Ans}(Q_0), \text{Ans}(Q^*) \mid Q^*) \\ = & \mathbf{E}_{\tilde{Q}}[H(\text{Ans}(Q_0) \mid Q^* = \tilde{Q}) + H(\text{Ans}(\tilde{Q}) \mid Q^* = \tilde{Q}) - H(\text{Ans}(Q_0), \text{Ans}(\tilde{Q}) \mid Q^* = \tilde{Q})] \end{aligned}$$

We can now apply Lemma 2 for any fixed \tilde{Q} and the event $\mathcal{E} = \{Q^* = \tilde{Q}\}$. Note that the density $\Pr[\mathcal{E}]$ is $2^{-k \cdot \Omega(\varepsilon \lg \frac{1}{\varepsilon})}$ which constant probability over the choice of \tilde{Q} . Thus, the lemma applies for small enough ε . We conclude that $H(\text{Ans}(Q_0) \mid \mathcal{E}) + H(\text{Ans}(\tilde{Q}) \mid \mathcal{E}) - H(\text{Ans}(Q_0), \text{Ans}(\tilde{Q}) \mid \mathcal{E}) = \Omega(k)$ with constant probability over \tilde{Q} . Thus, the expectation is also $\Omega(k)$.

Plugging our result into (1), the size of the encoding becomes $n + 3P + k \cdot O(\varepsilon \lg \frac{1}{\varepsilon}) - \Omega(k)$. Setting $k = \gamma P$ for a large constant γ , and ε a small enough constant, the negative $\Omega(k)$ term is double the positive terms. Thus, the encoding size is $n - \Omega(k)$, a contradiction.

References

- [CM96] David R. Clark and J. Ian Munro. Efficient suffix trees on secondary storage. In *Proc. 7th ACM/SIAM Symposium on Discrete Algorithms (SODA)*, pages 383–391, 1996.
- [GGG⁺07] Alexander Golynski, Roberto Grossi, Ankur Gupta, Rajeev Raman, and S. Srinivasa Rao. On the size of succinct indices. In *Proc. 15th European Symposium on Algorithms (ESA)*, pages 371–382, 2007.
- [GM03] Anna Gál and Peter Bro Miltersen. The cell probe complexity of succinct data structures. In *Proc. 30th International Colloquium on Automata, Languages and Programming (ICALP)*, pages 332–344, 2003.
- [Gol07] Alexander Golynski. Optimal lower bounds for rank and select indexes. *Theoretical Computer Science*, 387(3):348–359, 2007. See also ICALP’06.
- [Gol09] Alexander Golynski. Cell probe lower bounds for succinct data structures. In *Proc. 20th ACM/SIAM Symposium on Discrete Algorithms (SODA)*, pages 625–634, 2009.
- [GRR08] Alexander Golynski, Rajeev Raman, and S. Srinivasa Rao. On the redundancy of succinct data structures. In *Proc. 11th Scandinavian Workshop on Algorithm Theory (SWAT)*, 2008.
- [Jac89] Guy Jacobson. Space-efficient static trees and graphs. In *Proc. 30th IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 549–554, 1989.
- [Mil05] Peter Bro Miltersen. Lower bounds on the size of selection and rank indexes. In *Proc. 16th ACM/SIAM Symposium on Discrete Algorithms (SODA)*, pages 11–12, 2005.

- [MRR01] J. Ian Munro, Venkatesh Raman, and S. Srinivasa Rao. Space efficient suffix trees. *Journal of Algorithms*, 39(2):205–222, 2001. See also FSTTCS’98.
- [Mun96] J. Ian Munro. Tables. In *Proc. 16th Conference on the Foundations of Software Technology and Theoretical Computer Science (FSTTCS)*, pages 37–40, 1996.
- [Pag01] Rasmus Pagh. Low redundancy in static dictionaries with constant query time. *SIAM Journal on Computing*, 31(2):353–363, 2001. See also ICALP’99.
- [Păt08] Mihai Pătraşcu. Succincter. In *Proc. 49th IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 305–313, 2008.
- [RRR02] Rajeev Raman, Venkatesh Raman, and S. Srinivasa Rao. Succinct indexable dictionaries with applications to encoding k -ary trees and multisets. In *Proc. 13th ACM/SIAM Symposium on Discrete Algorithms (SODA)*, pages 233–242, 2002.